

# Regression

Ruwen Qin  
Stony Brook University

November 6, 2025

Data used in this module:

- Concrete\_Data.csv

Python notebook used in this module

- Regression.ipynb

Considering an entity characterized by an  $N$  dimensional feature representation:  $\mathbf{x} = [x_1, \dots, x_N]$ . We would like to infer measurements of interest,  $\mathbf{y} = [y_1, \dots, y_P]$ , from the features. Regression is one of the most widely used of all statistical methods to establish the mapping between  $\mathbf{x}$  and  $\mathbf{y}$ . In a regression model, the feature representation is the input, named predictors or independent variables. The measurements we would like to infer are the outputs and named responses or dependent variables. The responses and predictors can take both numerical and categorical values. This learning module is developed mainly based on references [1, 2].

## 1 Multiple Linear Regression

Multiple Linear Regression (MLR) handles the special case where there is one response and multiple predictors.

### 1.1 The Data and Hypothesis Set

Consider a training dataset with  $M$  data points. Each data point consists of  $N$  predictors and one response. The dataset in the matrix form includes

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_M \end{bmatrix} = \begin{bmatrix} x_{1,1} & \dots & x_{1,N} \\ \vdots & \ddots & \vdots \\ x_{M,1} & \dots & x_{M,N} \end{bmatrix}, \quad (1)$$

and

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_M \end{bmatrix}, \quad (2)$$

The multiple linear regression model relating the response  $y_i$  to predictors  $x_1, \dots, x_N$  is

$$y_i = w_0 + w_1 x_{i,1} + \dots + w_N x_{i,N} + \epsilon_i, \quad (3)$$

for  $i = 1, \dots, M$ .  $\epsilon_i$  in (3) is called the residual, noise, disturbance, or error.

Let

$$\mathbf{Z} = [\mathbf{1} \ \mathbf{X}] = \begin{bmatrix} \mathbf{z}_1 \\ \vdots \\ \mathbf{z}_M \end{bmatrix} = \begin{bmatrix} 1 & x_{1,1} & \dots & x_{1,N} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{M,1} & \dots & x_{M,N} \end{bmatrix}, \quad (4)$$

$$\boldsymbol{\epsilon} = \begin{bmatrix} \epsilon_1 \\ \vdots \\ \epsilon_M \end{bmatrix}, \quad (5)$$

$$\mathbf{w} = \begin{pmatrix} w_0 \\ \vdots \\ w_N \end{pmatrix}. \quad (6)$$

Then, the multiple linear regression model can be expressed as

$$\mathbf{y} = \mathbf{Z}\mathbf{w} + \boldsymbol{\epsilon}. \quad (7)$$

Assumptions of the linear regression model are:

1. Linearity of the conditional expectation:  $E(y_i | \mathbf{z}_i; \mathbf{w}) = \mathbf{z}_i \mathbf{w}$ ;
2. Independent noise:  $\epsilon_1, \dots, \epsilon_M$  are independent;
3. Constant variance:  $\text{Var}(\epsilon_i) = \sigma_\epsilon^2$  for all  $i$ ;
4. Gaussian noise:  $\epsilon_i$  is normally distributed for all  $i$ .

### 1.1.1 Least Squares Estimator (LSE)

The least-squares estimator (LSE) for  $\mathbf{w}$  minimizes the square of  $L_2$  distance between the true responses  $\mathbf{y}$  and their predictions  $\hat{\mathbf{y}}$ .

The LSE for  $\mathbf{w}$  is:

$$\hat{\mathbf{w}} = (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T \mathbf{y}. \quad (8)$$

Thereby, the vector of fitted response values is

$$\hat{\mathbf{y}} = \mathbf{Z} \hat{\mathbf{w}}, \quad (9)$$

.

### 1.1.2 Regression Residuals

Let

$$\hat{\boldsymbol{\epsilon}} = \mathbf{y} - \hat{\mathbf{y}}, \quad (10)$$

and  $\hat{\epsilon}_i$  denote the  $i$ th element of  $\hat{\boldsymbol{\epsilon}}$ . Then, an unbiased estimator of  $\sigma_{\epsilon}^2$  is

$$\hat{\sigma}_{\epsilon}^2 = \frac{\sum_{i=1}^M \hat{\epsilon}_i^2}{M - N - 1} = \frac{\sum_{i=1}^M (y_i - \hat{y}_i)^2}{M - N - 1}. \quad (11)$$

### 1.1.3 Confidence Interval Estimation of Regression Coefficients

The point estimator of  $\mathbf{w}$  in (8) is random. An interval estimation of the regression coefficients at the  $1 - \alpha$  confidence level can be obtained:

$$\hat{\mathbf{w}} \pm t_{1-\alpha/2, M-N-1} \text{se}(\hat{\mathbf{w}}) \quad (12)$$

where  $M - N - 1$  is the degrees of freedom for the regression residuals,  $t_{1-\alpha/2, M-N-1}$  is the t value providing an area of  $\alpha/2$  in the right tail of the t-distribution with  $M - N - 1$  degrees of freedom. If the interval estimate for a coefficient does not contains zero, the corresponding predictor has an impact on the response variable.

## 1.2 Analysis of Variance

### 1.2.1 Degree of Freedom

There are degrees of freedom (dof) associated with each of these sources of variation. The dof for regression is  $N$ , the number of predictors. The total degrees of freedom is  $M - 1$ . The degrees of freedom for regression residuals is  $M - N - 1$ .

### 1.2.2 Variance Partitioning

The total variation in  $\mathbf{y}$  (SST) can be partitioned into two parts: the variation that can be predicted by predictors (SSR) and the variation that cannot be predicted (SSE). Mathematically,

$$\begin{aligned} \text{SST} &= \sum_{i=1}^M (y_i - \bar{y})^2 \\ &= \sum_{i=1}^M (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^M (y_i - \hat{y}_i)^2 \\ &= \text{SSR} + \text{SSE} \end{aligned} \quad (13)$$

where  $\bar{y}$  is the sample mean computed based on  $\mathbf{y}$ .

### 1.2.3 Mean Squares of Error

The mean sum of squares is defined to be the ratio of sum of squares to its degrees of freedom.

Therefore, the mean squares total is

$$\text{MST} = \frac{\text{SST}}{M - 1}, \quad (14)$$

the mean squares of regression is

$$\text{MSR} = \frac{\text{SSR}}{N}, \quad (15)$$

and the mean squares of error is

$$\text{MSE} = \frac{\text{SSE}}{M - N - 1}. \quad (16)$$

which is the unbiased estimator of  $\sigma_\epsilon^2$ .

### 1.2.4 R-Squared, $R^2$

R-squared, denoted by  $R^2$ , is

$$R^2 = \frac{\text{SSR}}{\text{SST}} = 1 - \frac{\text{SSE}}{\text{SST}}, \quad (17)$$

and it measures the proportion of the total variation in  $\mathbf{y}$  that can be linearly predicted by predictors  $\mathbf{X}$ .

### 1.2.5 F Test

F test of goodness-of-fit uses the test statistic  $F$ :

$$F = \frac{\text{MSR}}{\text{MSE}} \quad (18)$$

to examine the hypotheses:

$$H_0: w_0 = \dots = w_N = 0; \forall j$$

$$H_1: \exists w_j \neq 0 \ j \in \{0, \dots, N\}$$

Rejecting the null hypothesis means a regression relationship exists between the predictors and the response variable. The p value corresponding to this test statistic is the probability that a random value drawn from the F distribution with  $N$  and  $M - N - 1$  degrees of freedom is greater than the test statistic. The larger the test statistic, the smaller the p value and so the chance of making a mistake in rejecting the null hypothesis.

## 1.3 Prediction Accuracy

The key concept associated with measuring forecast accuracy is forecast error. Various metrics are used to determine how well a particular forecasting method performs. Given  $\mathbf{x}_i$  from the test set, a regression model predicts the response  $\hat{y}_i$ . Multiple metrics are defined for assessing the prediction error for a numerical response:

Mean Absolute Error (MAE)

$$\text{MAE} = \frac{\sum_{i=1}^S |y_i - \hat{y}_i|}{S}, \quad (19)$$

Mean Square Error (MSE)

$$\text{MSE} = \frac{\sum_{i=1}^S (y_i - \hat{y}_i)^2}{S}, \quad (20)$$

Root Mean Square Error (RMSE)

$$\text{RMSE} = \sqrt{\frac{\sum_{i=1}^S (y_i - \hat{y}_i)^2}{S}}, \quad (21)$$

where  $S$  is the size of test dataset.

## 2 An Example of Multiple Linear Regression

### 2.1 The Problem

Concrete has 8 attributes:

1. Cement – quantitative – kg in a  $\text{m}^3$  mixture
2. Blast Furnace Slag – quantitative – kg in a  $\text{m}^3$  mixture
3. Fly Ash – quantitative – kg in a  $\text{m}^3$  mixture
4. Water – quantitative – kg in a  $\text{m}^3$  mixture
5. Superplasticizer – quantitative – kg in a  $\text{m}^3$  mixture
6. Coarse Aggregate – quantitative – kg in a  $\text{m}^3$  mixture
7. Fine Aggregate – quantitative – kg in a  $\text{m}^3$  mixture
8. Age – quantitative – Day (1~365)

We would like to build a regression model that either uses these hand-crafted features or extracted ones to predict concrete's compressive strength in megapascal (MPa).

### 2.2 Ordinary Linear Regression

#### 2.2.1 The Attained Model

We fit a linear regression model that predicts the compressive strength of concrete using the 8 principal components extracted from the eight concrete attributes as predictors. The 95% confidence intervals of the coefficients for C1 and C2 contain zero, indicating these two Cs may not be important predictors.

	Coefficient	Margin of error
1 $w_0$	35.666	0.722
2 $w_1$	-0.156	0.474
3 $w_2$	-0.468	0.605
4 $w_3$	7.375	0.636
5 $w_4$	1.329	0.719
6 $w_5$	-5.037	0.739
7 $w_6$	7.749	0.809
8 $w_7$	-7.370	1.680
9 $w_8$	11.752	4.150

#### 2.2.2 Residuals

First, we visually checked the residuals from the regression to make sure assumptions about residuals for the regression model hold. If needed, quantitative test results for the regression model can be calculated and they are more precise information.

Figure 1 is the plots of residuals attained on the training and testing datasets. The plots indicate residuals are independent.

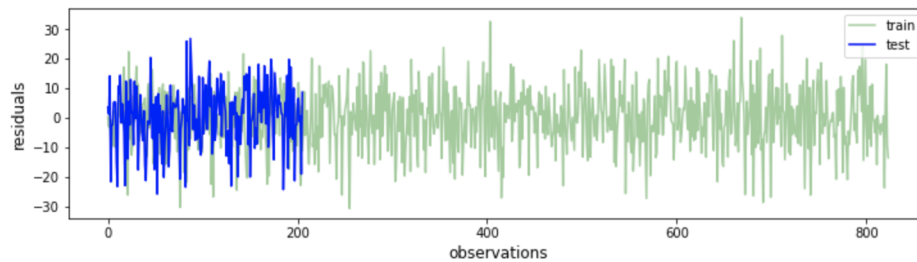
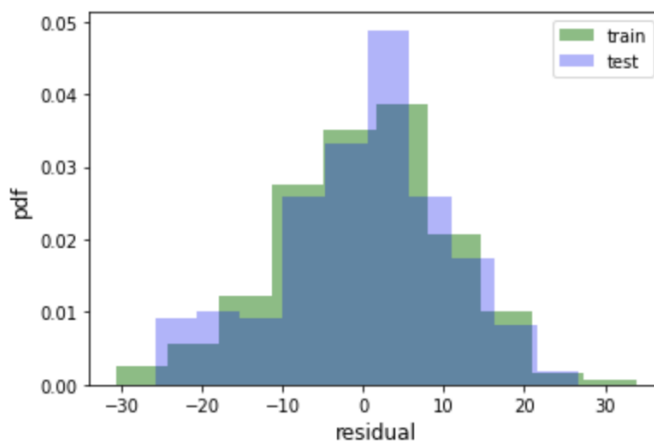


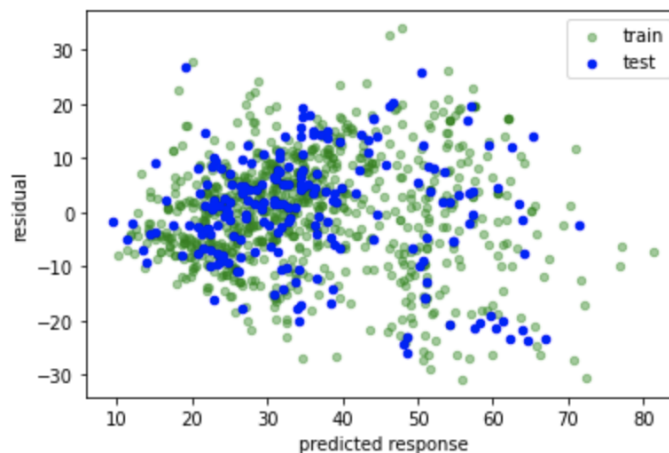
Figure 1: Residuals

Figure 2 shows the distributions of residuals obtained on the training and the testing datasets. The distributions are close to a normal distribution.



**Figure 2:** Residual Distribution

Figure 3 shows residuals at predicted response values. The plot shows the distribution of residuals increases when the predicted response value increases. Residuals are not homogeneously distributed. This raises an alarm. Therefore, the current linear model might not be the best choice. We will revisit this issue in later sections. In the following, we continued evaluating this linear model.



**Figure 3:** Residuals vs. predicted response

### 2.2.3 Goodness-of-Fit

In this example, the average  $R^2$  achieved on the training dataset using a five-fold cross validation is 56.41%, meaning the amount of data variation captured by the model.  $R^2_{adj}$  is 55.97%, slightly smaller than  $R^2$ , indicating the model is not that complex.

```
1 AIC:6,051
2 BIC:6,093
3 The average R squared value for training: 56.41%
4 The average adjusted R squared value for training: 55.97%
```

In testing the model on a test dataset, the accuracy is measured as below:

```
1 MAE: 7.618
2 MSE: 99.141
```

```

3 RMSE: 9.957
4 R Square: 0.616

```

### 3 Polynomial Linear Regression

Polynomial linear regression fits a polynomial function of predictors  $\mathbf{x} = [x_1, \dots, x_N]$  to predict the response  $y$ .

For example, a polynomial in three predictors  $[x_1, x_2, x_3]$  with degree of 2 has 10 ( $=C_2^{3+2}$ ) polynomial terms: 1,  $x_1$ ,  $x_2$ ,  $x_3$ ,  $x_1x_2$ ,  $x_1x_3$ ,  $x_2x_3$ ,  $x_1^2$ ,  $x_2^2$ ,  $x_3^2$ .

We can fit a polynomial regression function using the linear regression model.

### 4 Binary Logistic Regression

#### 4.1 The Data and Hypothesis Set

A training dataset contains  $M$  data points. Each data point consists of the observations of  $N$  predictors and a binary response. Binary logistic regression models the conditional probability that a binary response is 1 given predictors and their coefficients. The binary logistic regression can take the following form:

$$p_i = p(y_i = 1 | \mathbf{z}_i; \mathbf{w}) = \text{sigmoid}(\mathbf{z}_i \mathbf{w}) = \frac{1}{1 + \exp(-\mathbf{z}_i \mathbf{w})}, \quad (22)$$

where  $\mathbf{w}$  is the vector of regression coefficients

$$\mathbf{w} = \begin{bmatrix} w_0 \\ \vdots \\ w_N \end{bmatrix}, \quad (23)$$

and  $\mathbf{z}_i = [1 \ \mathbf{x}_i]$ . We thereby define  $\mathbf{Z}$  by appending a column of ones to the left of  $\mathbf{X}$ :

$$\mathbf{Z} = [\mathbf{1} \ \mathbf{X}] = \begin{bmatrix} \mathbf{z}_1 \\ \vdots \\ \mathbf{z}_M \end{bmatrix} = \begin{bmatrix} 1 & x_{1,1} & \dots & x_{1,N} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{M,1} & \dots & x_{M,N} \end{bmatrix}. \quad (24)$$

By learning from the training dataset, the regression coefficients  $\mathbf{w}$  can be estimated, for example, using maximum likelihood estimation.

#### 4.2 An Example

Logistic regression is commonly used in classification problems. Let's consider a problem similar to the concrete example we studied. The response  $y$  is a binary variable indicates the class of a concrete slab: "defective" vs. "non-defective".  $y$  takes the value 1 meaning defective, and 0 meaning non-defective. We split the data by 80-20. That is, 80% of the data are used for training a model and 20% for testing the developed model.

The maximum likelihood estimates of regression coefficients are:

```

1 Coefficient
2 w0 -3.253
3 w1 -1.644
4 w2 0.924
5 w3 -1.609
6 w4 2.284
7 w5 2.062
8 w6 -5.268
9 w7 1.082
10 w8 -2.021

```

The predicted class is

$$\hat{y}_i = \begin{cases} 1 & \text{if } p_i > 0.5 \\ 0 & \text{if } p_i \leq 0.5. \end{cases} \quad (25)$$

Figure 4 is the confusion matrix. Among the 136 true zeros, 123 are predicted correctly, and 13 are predicted mistakenly. Among the 65 true ones, 57 are predicted correctly, and 8 are predictably mistakenly. The accuracy is 89.55% (= (123 + 57)/(123 + 13 + 9 + 57)).

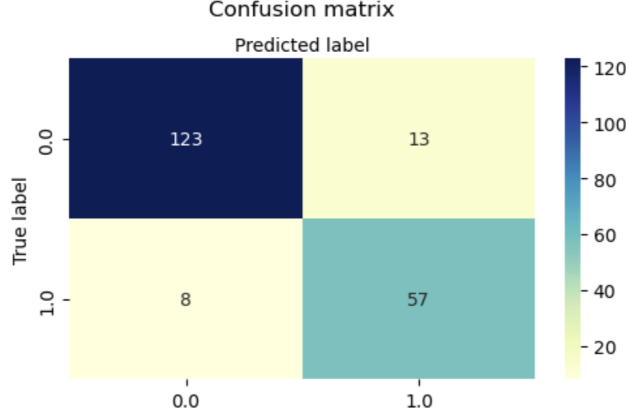


Figure 4: Confusion matrix

## 5 Closing Remark

In this learning module, we primarily introduced several popularly used regression models falling in the category of *Generalized Linear Models* (GLM). A GLM maps a linear combination of predictors  $\mathbf{x}_i \mathbf{w}$  to the expected value of response  $E[y_i | \mathbf{x}_i]$  through a link function  $l(\mathbf{x}_i \mathbf{w}) = E[y_i | \mathbf{x}_i; \mathbf{w}]$ .

We considered various distributions of the response in regression: the response in the multiple linear regression follows a univariate Gaussian distribution, and the Bernoulli distribution in the binary logistic regression. There are a few more distributions of responses that this learning module did not cover, such as responses following binomial distribution or Poisson distribution. Regression models for those can be developed similarly.

## References

- [1] Kevin P. Murphy. *Probabilistic Machine Learning: An introduction*. MIT Press, 2022.
- [2] Kevin P. Murphy. *Probabilistic Machine Learning: Advanced Topics*. MIT Press, 2023.